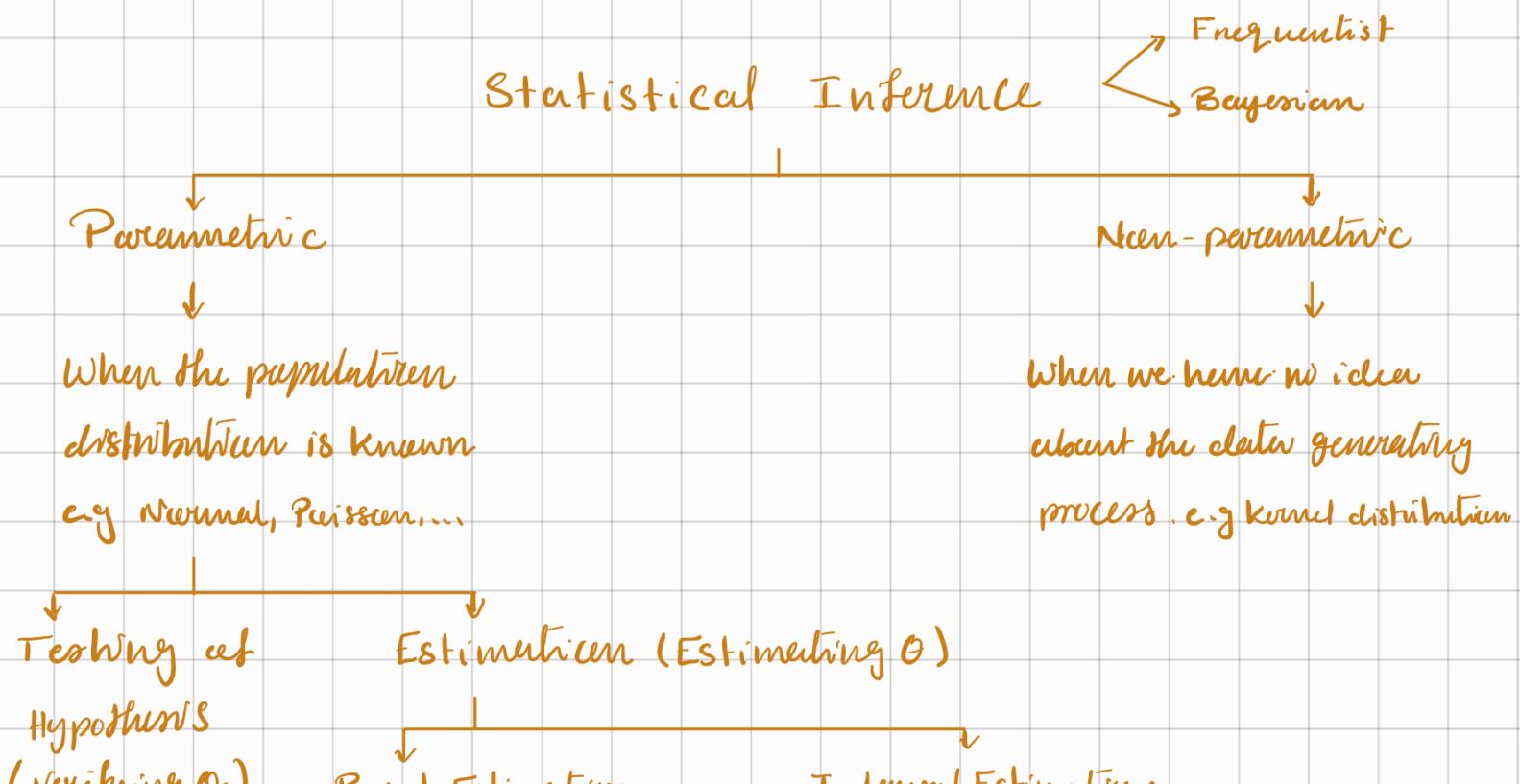


- Statistical Inference refers to drawing conclusion based on evidence obtained from the data .
- Statistical data consists of variability and uncertainty statistics models variability and quantifies uncertainty.



\* Main challenges :

- 1) How to summarize the information in the data using formal mathematical tools ? Min Q1 Median Q3 Max , Mean
- 2) How to use these summaries to answer questions about the phenomenon at interest ?

- \* Examples :
- 1) Survival analysis
  - 2) Biostatistics
  - 3) NASA dataset produced in their software defect detection trials
  - 4) Daily price of the saleable stocks for a company

\* Data :  $n_1, \dots, n_n$  where  $n_j$  can either be a scalar or a vector.  
In stat inference, this sample is interpreted as a realization of RVs.

\* Notation :  $\tilde{X} = (x_1, \dots, x_n)^T$  is a vector of RVs  
 $\tilde{n} = (n_1, \dots, n_n)^T$  is the sample of observations

\* Statistic : Function of sample that helps in data summarization (reduction)  
 $T(\tilde{n}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ;  $1 \leq m \leq n$  e.g.  $m=2$  ( $\mu, \sigma^2$ ) for normal  
 Instead of reporting the entire sample, only report the value of statistic  
 Visual 5-point summary statistic : Boxplot

\* Point Estimation : Uses sample data to calculate a single value (statistic) which serves as the best estimate of the unknown (fixed / random) population parameter
 

- $(x_1, \dots, x_n) \sim F_\theta$ ,  $\theta \in \Omega$ , Guess  $\theta$ ? or Guess  $\sigma(\theta)$ ?
- $T(\tilde{n})$  is called the point estimate of  $\theta$ .

- \* Interval Estimation : Interval of possible (probable) values of an unknown population parameter.  $L(\hat{\theta}) \leq \theta \leq U(\hat{\theta})$   
 $\theta \in [L(\hat{\theta}), U(\hat{\theta})]$  is called an interval estimator.  
 Bayesian : Credible interval (deal with prior info)  
 Frequentist : Confidence interval / Prediction Intervals

\* Parametric Models :  $\{f(n|\theta) \mid \theta \in \Omega\}$  with  $k$  parameters  $\Omega \subseteq \mathbb{R}^k$  is the parameter space to which the parameters belong and  $f(n|\theta)$  is the PDF / PMF for the distribution ( $k$  is finite)

\* Point Estimation : 1) Method of Moments  
 2) Maximum Likelihood Estimation

\* Measure of Goodness : 1) Unbiasedness      3) Efficiency  
 2) Consistency      4) Completeness

\* Method of Moments :  $x_1, \dots, x_n \stackrel{iid}{\sim} f(n|\theta)$ . How to estimate  $\theta$ ?

$X \sim f(n|\theta)$ ;  $\theta \in \mathbb{R}^k$  so we consider the first  $k$  moments of the distribution of  $X$ . To estimate  $\theta$ , we use observed sample moments

$$\begin{aligned} p_1 &= E(X) & \hat{p}_1 &= \frac{1}{n} (x_1 + \dots + x_n) = \bar{x} && \text{equals the sample moments to the theoretical moments} \\ p_2 &= E(X^2) & \hat{p}_2 &= \frac{1}{n} (x_1^2 + \dots + x_n^2) \\ \vdots & & \vdots & & & \\ p_k &= E(X^k) & \hat{p}_k &= \frac{1}{n} (x_1^k + \dots + x_n^k) \end{aligned}$$

$$\hat{p}_2 - \hat{p}_1^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$$

1)  $X \sim \text{Pois}(\lambda)$

$$\mathbb{E}(X) = \lambda \Rightarrow \hat{\lambda} = \frac{1}{n} (x_1 + \dots + x_n) = \bar{x}$$

2)  $X \sim \exp(\lambda)$

$$\mathbb{E}(X) = \frac{1}{\lambda} \Rightarrow \frac{1}{\hat{\lambda}} = \frac{1}{n} (x_1 + \dots + x_n) \Rightarrow \hat{\lambda} = \frac{1}{\bar{x}}$$

3)  $X \sim N(\mu, \sigma^2)$

$$\mathbb{E}(X) = \mu$$

$$\mathbb{E}(X^2) = \text{Var}(X) + \mathbb{E}(X)^2 = \sigma^2 + \mu^2$$

$$\hat{\mu}_1 = \mu$$

$$\Rightarrow \hat{\mu} = \bar{x}$$

$$\hat{\mu}_2 = \sigma^2 + \mu^2 = \sigma^2 + \hat{\mu}_1^2$$

$$\Rightarrow \hat{\sigma}^2 = \hat{\mu}_2 - \hat{\mu}_1^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$$

$$= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

4)  $X \sim \text{Gamma}(\alpha, \beta)$

$$\mathbb{E}(X) = \frac{\alpha}{\beta} = \hat{\mu}_1$$

$$\mathbb{E}(X^2) = \text{Var}(X) + \mathbb{E}(X)^2 = \frac{\alpha}{\beta^2} + \frac{\alpha^2}{\beta^2} = \frac{\alpha + \alpha^2}{\beta^2} = \hat{\mu}_2$$

$$\frac{\hat{\alpha}}{\hat{\beta}} = \bar{x} \Rightarrow \hat{\alpha} = \hat{\beta} \bar{x} = \frac{n \bar{x}^2}{\sum (x_i - \bar{x})^2}$$

$$\begin{aligned} \frac{\hat{\alpha} + \hat{\alpha}^2}{\hat{\beta}^2} &= \frac{1}{n} \sum_{i=1}^n x_i^2 \Rightarrow \hat{\beta} \bar{x} + \frac{\hat{\alpha}^2}{\hat{\beta}^2} \bar{x}^2 = \frac{1}{n} \sum x_i^2 \\ &\Rightarrow \bar{x} \left( \frac{1}{\hat{\beta}} + \bar{x} \right) = \frac{1}{n} \sum x_i^2 \\ &\Rightarrow \frac{1}{\hat{\beta}} + \bar{x} = \frac{1}{n \bar{x}} \sum x_i^2 \end{aligned}$$

$$\Rightarrow \frac{1}{\hat{\beta}} = \frac{1}{n \bar{x}} \sum x_i^2 - \bar{x} = \frac{1}{\bar{x}} \left( \frac{1}{n} \sum x_i^2 - \bar{x}^2 \right)$$

$$\Rightarrow \hat{\beta} = \frac{n \bar{x}}{\sum (x_i - \bar{x})^2}$$

\* Method of MLE :  $L(\boldsymbol{\theta} | \tilde{x}) = \prod_{i=1}^n f(x_i | \boldsymbol{\theta})$  likelihood function

MLE is finding  $\hat{\theta}$  that maximizes  $L(\theta|x)$  ( $L(\theta|x)$  measures how likely is a particular value of  $\theta$  given  $\tilde{x}$  is observed)

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmax}} (L(\theta|x)) \quad \text{Optimization Problem}$$

We can maximize  $l(\theta) = \log(L(\theta|x)) = \sum_{i=1}^n \log(f(x_i|\theta))$  instead

i)  $X \sim \text{Bin}(p)$

$$\begin{aligned} f(p|x) &= \prod_{i=1}^n f(x_i|p) \\ &= \prod_{i=1}^n p^{n_i} (1-p)^{1-n_i} \\ &= p^{\sum n_i} (1-p)^{n - \sum n_i} \end{aligned}$$

$$l(p) = \sum_{i=1}^n n_i \log(p) + (n - \sum_{i=1}^n n_i) \log(1-p)$$

$$\frac{\partial l(p)}{\partial p} = \frac{1}{p} \sum_{i=1}^n n_i - \frac{1}{1-p} (n - \sum_{i=1}^n n_i) \stackrel{\text{set}}{=} 0$$

$$(1-p) \sum n_i - p(n - \sum_{i=1}^n n_i) = 0$$

$$\sum_{i=1}^n n_i - pn = 0$$

$$\hat{p}_{MLE} = \frac{1}{n} \sum_{i=1}^n n_i = \bar{x}$$

$$\text{Checking: } \frac{\partial^2}{\partial p^2} l(p|x) = \frac{\partial}{\partial p} \left( \frac{\sum n_i - np}{p(1-p)} \right) < 0 \Rightarrow$$

$\hat{p} = \bar{x}$  is the MLE at  $p$

2) Let  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$

$$\ell(\mu, \sigma^2 | X_i) = \prod_{i=1}^n f(X_i | \mu, \sigma^2)$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(X_i - \mu)^2}{2\sigma^2}\right)$$

$$= (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\sum_{i=1}^n \frac{(X_i - \mu)^2}{2\sigma^2}\right)$$

$$\ell(\mu, \sigma^2) = \frac{-n}{2} \log(2\pi) - \frac{n}{2} \log \sigma^2 - \sum_{i=1}^n \frac{(X_i - \mu)^2}{2\sigma^2}$$

$$\frac{\partial \ell(\mu, \sigma^2)}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu) = \frac{1}{\sigma^2} (\sum_{i=1}^n X_i - n\mu) \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \sum_{i=1}^n X_i = n\mu \Rightarrow \hat{\mu}_{MLE} = \bar{X}$$

$$\frac{\partial \ell(\mu, \sigma^2)}{\partial \sigma^2} = \frac{-n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (X_i - \hat{\mu}_{MLE})^2 \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \sigma^2_{MLE} = \sum_{i=1}^n (X_i - \bar{X})^2$$

$$\Rightarrow \sigma^2_{MLE} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

3)  $X_1, \dots, X_n \sim P(\lambda)$

$$\ell(\lambda | X) = \prod_{i=1}^n f(X_i | \lambda)$$

$$= \prod_{i=1}^n \frac{e^{-\lambda} \lambda^{X_i}}{X_i!}$$

$$= \frac{e^{-n\lambda} \lambda^{\sum X_i}}{\prod_{i=1}^n (X_i)!}$$

$$\ell(\lambda) = -n\lambda + \log(\lambda) \sum_{i=1}^n X_i - \sum_{i=1}^n \log(X_i!)$$

$$\frac{\partial \ell(\lambda)}{\partial \lambda} = -n + \frac{1}{\lambda} \sum_{i=1}^n x_i \stackrel{\text{set}}{=} 0 \Rightarrow \hat{\lambda}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$$

4) Exponential :  $f(x_i | \mu, \lambda) = \lambda e^{-\lambda(x_i - \mu)} = \lambda e^{-\lambda x_i + \lambda \mu}, n \geq 1$   
 $\mu \in \mathbb{R}$

$$\begin{aligned} L(x_i | \mu, \lambda) &= \prod_{i=1}^n f(x_i | \mu, \lambda) \\ &= \prod_{i=1}^n \lambda e^{-\lambda(x_i - \mu)} \\ &= \lambda^n e^{-\lambda \sum_{i=1}^n (x_i - \mu)} \end{aligned}$$

$$\ell(\mu, \lambda) = n \log \lambda - \lambda \sum_{i=1}^n (x_i - \mu)$$

- As  $e^y$  is an increasing function then  $\forall y, y' ; y < y' \quad e^y < e^{y'}$  so to maximize such a function we need to consider  $\hat{\mu}_{MLE} = \min_{1 \leq i \leq n} (x_i)$

- $\frac{\partial \ell(\mu, \lambda)}{\partial \lambda} = \frac{n}{\lambda} - \sum_{i=1}^n (x_i - \hat{\mu}) \stackrel{\text{set}}{=} 0 \Rightarrow$  as  $\mu \leq x_i \forall i$

$$n = \lambda \sum_{i=1}^n (x_i - \hat{\mu}) \Rightarrow$$

$$\hat{\lambda}_{MLE} = \frac{n}{\sum_{i=1}^n x_i - n \hat{x}_{(1)}} = \frac{1}{\frac{1}{n} \sum_{i=1}^n x_i - \hat{x}_{(1)}} = \frac{1}{\bar{x} - \hat{x}_{(1)}}$$

5)  $X \sim \text{Uniform}[a, b]$   $f_{x_i} = \frac{1}{b-a} \mathbf{1}(a \leq x_i \leq b)$

$$\begin{aligned} L(a, b | X) &= \prod_{i=1}^n f(x_i | a, b) \\ &= \frac{1}{(b-a)^n} \mathbf{1}(a \leq x_{(1)} \leq x_{(n)} \leq b) \\ &= \frac{1}{(b-a)^n} \mathbf{1}(a \leq x_{(1)}) \mathbf{1}(x_{(n)} \leq b) \end{aligned}$$

As  $n \mapsto \frac{1}{n}$  is a decreasing function,  $\hat{a}_{MLE} = x_{(1)}$  and  $\hat{b}_{MLE} = x_{(n)}$

(6)  $X_1, \dots, X_k \sim \text{Multinomial}(n, (p_1, \dots, p_k))$  (Important for final)

$$f(x_1, \dots, x_k | p_1, \dots, p_k) = \binom{n}{x_1, \dots, x_k} p_1^{x_1} \cdots p_k^{x_k}, \quad \sum_{i=1}^k p_i^{x_i} = 1 \text{ and}$$

$$\ell(p_1, \dots, p_k) = \log \left( \binom{n}{x_1, \dots, x_k} \right) + \sum_{i=1}^k x_i \log(p_i) \quad \sum_{i=1}^k n_i = n$$

$$\begin{aligned} F(p_1, \dots, p_k) &= \ell(p_1, \dots, p_k) + \lambda g(p_1, \dots, p_k) \\ &= \log \left( \binom{n}{x_1, \dots, x_k} \right) + \sum_{i=1}^k x_i \log(p_i) + \lambda \left( \sum_{i=1}^k p_i - 1 \right) \end{aligned}$$

$$\frac{\partial F(p_1, \dots, p_k)}{\partial p_i} = \frac{x_i}{p_i} + \lambda = 0 \Rightarrow \hat{p}_i = \frac{-x_i}{\lambda} \quad \forall i \in \{1, \dots, k\}$$

$$\text{Choose } \lambda = -n \text{ then } \forall i \in \{1, \dots, k\} \quad \hat{p}_i = \frac{x_i}{n} \text{ so } \sum_{i=1}^k \frac{x_i}{n} = \frac{n}{n} = 1$$

thus this  $\lambda$  verifies the condition.

$$\text{Alternatively, } \sum_{i=1}^k \hat{p}_i = \sum_{i=1}^k \frac{-x_i}{\lambda} = \frac{-n}{\lambda} = 1 \Rightarrow \lambda = -n$$

### \* The Tank Problem:

The enemy has an unknown number  $N$  of tanks which he has arbitrarily numbered  $1, 2, \dots, N$ . Spies have reported sighting 8 tanks with numbers

$$\tilde{n} = (137, 24, 86, 33, 92, 129, 17, 111)^T$$

Assume that sightings are independent and that each of the  $N$  tanks has a probability  $\frac{1}{N}$  of being observed at each sighting. What is the MLE of  $N$ ?

$$L(N | \tilde{n}) = \left\{ \frac{1}{N^8}, \quad N \geq 137 \right.$$

$$\hat{N} = x_{(8)} = 137$$

\* Why MLE ? 1) The estimators are consistent and asymptotically normal i.e for likelihood  $L(\theta | n)$  & n sample size

$$\hat{\theta}_{MLE} \xrightarrow{} \theta \text{ as } n \rightarrow +\infty$$

$$\sqrt{n} (\hat{\theta}_{MLE} - \theta) \xrightarrow{d} N(0, \Sigma^*)$$

$\Sigma^*$  is an estimate matrix called Fisher information matrix

So if we use MLE method, we can construct confidence

interval around  $\hat{\theta}_{MLE}$

2) Sometimes the moments do not exist so we cannot use the method at moments.

\* Sometimes we don't find a closed form solution at MLE

$$L(\alpha, \beta | X) = \prod_{i=1}^n f(x_i | \alpha, \beta)$$

$$= \prod_{i=1}^n \frac{\beta^\alpha}{\Gamma(\alpha)} x_i^{\alpha-1} e^{-\beta x_i}$$

$$= \frac{\beta^{n\alpha}}{\Gamma(\alpha)^n} \left( \prod_{i=1}^n x_i^{\alpha-1} \right) e^{-\beta \sum_{i=1}^n x_i}$$

$$l(\alpha, \beta) = n\alpha \log(\beta) - n \log(\Gamma(\alpha)) + (\alpha-1) \sum_{i=1}^n \log(x_i) - \beta \sum_{i=1}^n x_i$$

$$\frac{\partial l(\alpha, \beta)}{\partial \beta} = \frac{n\alpha}{\beta} - \sum_{i=1}^n x_i = 0 \Rightarrow n\alpha = \beta \sum_{i=1}^n x_i \\ \Rightarrow \hat{\beta}_{MLE} = \frac{\hat{\alpha}}{\bar{x}}$$

$$\frac{\partial l(\alpha, \beta)}{\partial \alpha} = n \log(\beta) - n \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} + \sum_{i=1}^n \log(x_i)$$

$$= n \log\left(\frac{\hat{\alpha}}{\bar{x}}\right) - n \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} + \sum_{i=1}^n \log(x_i)$$

- \* We will study optimization methods :
  - Newton Raphson (NR)
  - Gradient Descent (GD)

## i) Newton Raphson :

Taylor Expansion :  $f(n) = f(\alpha) + (n-\alpha) f'(\alpha) = 0 \Rightarrow$   
 $-f(\alpha) = (n-\alpha) f'(\alpha) \Rightarrow$   
 $n = \alpha - \frac{f(\alpha)}{f'(\alpha)}$

We need to start with an initial guess ( $\alpha^{(0)}$ ), which may be the method of moments estimator  $\alpha^{(0)} = \frac{n\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$

Then, we compute  $\alpha^{(1)}, \alpha^{(2)}, \dots, \alpha^{(n)}$  selection criterion is whether  $\alpha^{(i)}$  is better than  $\alpha^{(j)}$ .

Having computed  $\alpha^{(t)}$ ,  $t \in \{0, 1, \dots, n\}$ , we compute the next iteration  $\alpha^{t+1}$  by approximating the function with a linear equation using the 1<sup>st</sup> order Taylor Expansion around  $\hat{\alpha} = \alpha^{(t)}$ , and set  $\alpha^{(t+1)}$  as the value at  $\hat{\alpha}$  that solves this linear equation.

$$0 = \text{log}(\hat{\alpha}) - \frac{f'(\alpha)}{f(\alpha)} - \text{log}(\bar{x}) + \frac{1}{n} \sum_{i=1}^n \text{log} x_i$$

we have the same

$$f(\alpha) = \text{log}(\alpha) - \frac{f'(\alpha)}{f(\alpha)}$$

$$\alpha_{t+1} = \alpha_t - \frac{f'(\alpha_t)}{f''(\alpha_t)}$$

result as we originally

considered  $\frac{\partial \text{LL}(a, b)}{\partial a}$  and

then applied NR

$$\alpha^{(1)} = \alpha^{(0)} - \frac{f(\alpha_0)}{f'(\alpha_0)}$$

if we think  $\alpha^{(1)}$  is better then we update it

$$\alpha^{(2)} = \alpha^{(1)} - \frac{f(\alpha_1)}{f'(\alpha_1)}$$

## ii) Gradient Ascent :

Taylor Expansion :  $f(\theta) = f(\theta_0) + (\theta - \theta_0) f'(\theta_0) + \frac{(\theta - \theta_0)^2}{2} f''(\theta_0)$

Suppose  $f''(\theta_0)$  is -ve constant :  $f(\theta) = f(\theta_0) + (\theta - \theta_0) f'(\theta_0) - \frac{(\theta - \theta_0)^2}{2}$   
 Hence, we get  $\theta = \theta_0 + t f'(\theta_0)$  by setting the derivative to zero  
 $\downarrow$   
 Learning Rate

The choice of the learning rate is the most important consideration.

$$f'(\theta) = f'(\theta_0) - \frac{(\theta - \theta_0)}{t} = 0 \quad \text{as } f''(\theta_0) = \frac{-1}{\delta}$$

$$\Rightarrow \theta = \theta_0 + t f'(\theta_0)$$

$\downarrow$   
Learning Parameter

Concave Function : Reach global maxima

General Function : Reach local maxima

\* Higher Dimensions :  $\theta_{k+1} = \theta_k + t \nabla f(\theta_k)$

\* Bias, Variance and Mean Squared Error :

1) Single Parameter :  $x_1, \dots, x_n \stackrel{iid}{\sim} f_{\theta}(x)$ . We estimate  $\theta$  using MOM and MLE ( $\hat{\theta}$ ). How to determine if this estimate is good or bad?

$$\text{Bias} = E_{\theta}(\hat{\theta}) - \theta$$

$$\text{Standard Error} : \sqrt{\text{Var}_{\theta}(\hat{\theta})}$$

$$\text{MSE} = E_{\theta}((\theta - \hat{\theta})^2)$$

Bias measures how close the average value at  $\hat{\theta}$  to the true parameter  $\theta$   
 SE measures how variable is  $\hat{\theta}$  around this average value

$$\begin{aligned} \text{MSE} &= E_{\theta}((\theta - \hat{\theta})^2) \\ &= E_{\theta}(\theta^2 - 2\theta\hat{\theta} + \hat{\theta}^2) \\ &= E(\theta^2) - 2E(\theta)\hat{\theta} + E(\hat{\theta}^2) \end{aligned}$$

$$\begin{aligned}
 &= (\mathbb{E}(\hat{\alpha}^2) - 2\mathbb{E}(\hat{\alpha})\mathbb{E}(\hat{\alpha}) + \mathbb{E}(\hat{\alpha})^2) + \text{Var}(\hat{\alpha}) \\
 &= \text{Var}(\hat{\alpha}) + (\mathbb{E}(\hat{\alpha}) - \alpha)^2 \\
 &= \text{Var}(\hat{\alpha}) + \text{Bias}^2
 \end{aligned}$$

\* Consistency : An estimator  $\hat{\alpha}$  based on  $x_1, \dots, x_n$  is consistent if  
 $\hat{\alpha} \xrightarrow{P} \alpha$  as  $n \rightarrow +\infty$  ( $P(|\hat{\alpha} - \alpha| < \varepsilon) \rightarrow 1 \forall \varepsilon$ )

Working conditions :  $\mathbb{E}(\hat{\alpha}) = \alpha$   
 $\text{Var}(\hat{\alpha}) \xrightarrow{n \rightarrow +\infty} 0$

\* Examples :

1)  $\hat{\lambda} = \bar{x}$  by MOM and MLE ,  $x_1, \dots, x_n \sim S(\lambda)$

$$\mathbb{E}(\hat{\lambda}) = \mathbb{E}(\bar{x}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(x_i) = \lambda \quad \text{Unbiased}$$

$$\text{Var}(\hat{\lambda}) = \text{Var}(\bar{x}) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(x_i) = \frac{1}{n} \xrightarrow[n \rightarrow +\infty]{\downarrow} 0, \quad \text{sd}(\hat{\lambda}) = \sqrt{\frac{1}{n}} \quad \text{consistent}$$

$$\text{MSE} = \text{Bias}^2 + \text{Var}(\hat{\lambda}) = 0 + \frac{1}{n} = \frac{1}{n}$$

2)  $\hat{\lambda} = \frac{1}{\bar{x}}$  for MOM and MLE .  $x_1, \dots, x_n \sim \text{Exp}(\lambda)$

$$\text{Bias} : \mathbb{E}(\hat{\lambda}) = \mathbb{E}\left(\frac{n}{\sum x_i}\right) = \frac{n}{n-1} \quad \bar{x} \sim \text{Gamma}(n, n\lambda)$$

Jensen's Inequality :  $\mathbb{E}(\ell(\bar{x})) \leq \ell(\mathbb{E}(\bar{x}))$

$$\lambda = \frac{1}{\mathbb{E}(\bar{x})} \leq \mathbb{E}\left(\frac{1}{\bar{x}}\right) = \mathbb{E}(\hat{\lambda}) \text{ so it is unbiased}$$

3) Nonsense Unbiased Estimator

Let  $x \sim S(\lambda)$  ,  $\lambda > 0$  . We want to estimate the parameter  $\alpha = e^{-\beta\lambda}$

based on a sample of size 1.  $T(X) = (-\lambda)^X$

$$\begin{aligned} E(T(X)) &= E((- \lambda)^X) = \sum_{k \geq 0} \frac{(-\lambda)^k e^{-\lambda} \cdot \lambda^k}{k!} \\ &= e^{-\lambda} \sum_{k \geq 0} \frac{(-\lambda \lambda)^k}{k!} \\ &= e^{-\lambda} \cdot e^{-\lambda \lambda} \\ &= e^{-2\lambda} \end{aligned}$$

Hence, this is not a good estimator as when  $X$  is odd the estimate is negative but  $\lambda$  is a positive quantity.

\* **Remark K :** When  $n$  is large, asymptotic theory provides us with a more complete picture of the "accuracy" of  $\lambda$ . By CLT  $\sqrt{n}(\bar{\lambda} - \lambda) \xrightarrow{a} N(0, 1)$  as  $n \rightarrow +\infty$   
a : asymptotically

\* **Asymptotically Normal :**  $\sqrt{n}(\hat{\lambda} - \lambda) \xrightarrow{d} N(0, \sigma)$

Finding Asymptotic Variance :

$$f(n|\lambda) = \frac{e^{-\lambda} \lambda^n}{n!} \quad \text{log}(f(n|\lambda)) = -\lambda + n \text{log}(\lambda) - \text{log}(n!)$$

$$z(n, \lambda) = \frac{\partial}{\partial \lambda} (\text{log}(f(n|\lambda))) = \frac{n}{\lambda} - 1$$

$$z'(n, \lambda) = \frac{-n}{\lambda^2}$$

Fisher Information :  $I(\lambda) = -E_{\lambda}(z'(X, \lambda)) = -E\left(\frac{-X}{\lambda^2}\right) = \frac{1}{\lambda}$

$$\text{Hence, } \hat{\lambda} = \frac{a}{n} \sim N\left(\frac{a}{n}, \frac{1}{n^2}\right) \sim N\left(\lambda, \frac{1}{n}\right)$$

Hence,  $\ln(\frac{f(x|\theta)}{f(x|\theta_0)}) \sim N(0, \frac{1}{n I(\theta_0)})$

### \* Theorem :

Let  $\{f(u|\theta) | \theta \in \Omega\}$  be a parametric model  $\theta \in \mathbb{R}$  is a single parameter. Let  $x_1, \dots, x_n \stackrel{iid}{\sim} f(u|\theta_0)$  for  $\theta_0 \in \Omega$ . Let  $\hat{\theta}$  be the MLE based on  $x_1, \dots, x_n$ . Suppose certain regularity conditions hold.

- 1) All the pdf / pmf  $f(u|\theta)$  in the model have the same support.
- 2)  $\theta_0$  is an interior point (not on the boundary)
- 3) The log likelihood  $\ell(\theta)$  is differentiable w.r.t  $\theta$
- 4)  $\hat{\theta}$  is the unique value (MLE) at  $\theta \in \Omega$  that satisfies  $\ell'(\theta) = 0$

Then,  $\hat{\theta}$  is consistent and asymptotically normal, with

$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} N\left(0, \frac{1}{I(\theta_0)}\right)$$

$I(\theta) := V_\theta[z(x, \theta)] = -\mathbb{E}(z'(x, \theta))$  where  $z(x, \theta)$  is the score function.

### \* Proof :

$$\begin{aligned} \mathbb{E}(z(x, \theta)) &= \int z(x, \theta) f(x|\theta) dx \\ &= \int \frac{\partial}{\partial \theta} \log(f(x|\theta)) f(x|\theta) dx \\ &= \int \frac{\frac{\partial \ln f(x|\theta)}{\partial \theta}}{f(x|\theta)} \cdot f(x|\theta) dx \\ &= \int \frac{\partial}{\partial \theta} f(x|\theta) dx \end{aligned}$$

$$= \frac{\partial}{\partial \alpha} \int f(x|\alpha) dn \quad 1$$

$$= 0$$

$$\mathbb{E}(z'(x, \alpha)) = \int z'(x, \alpha) f(x|\alpha) dn$$

$$= \int \frac{\partial^2}{\partial^2 \alpha} \log(f(x|\alpha)) f(x|\alpha) dn$$

$$= \int \frac{\partial}{\partial \alpha} z(x, \alpha) f(x|\alpha) dn$$

$$= \frac{\partial}{\partial \alpha} \int z(x, \alpha) f(x|\alpha) dn$$

$$= \frac{\partial}{\partial \alpha} \mathbb{E}(z(x, \alpha))$$

$$\mathbb{E}(z'(x, \alpha)) = \frac{\partial}{\partial \alpha} \mathbb{E}(z(x, \alpha))$$

$$0 = \frac{\partial}{\partial \alpha} \int z(x, \alpha) f(x|\alpha) dn$$

$$= \int \frac{\partial}{\partial \alpha} (z(x, \alpha) f(x|\alpha)) dn$$

$$= \int z'(x, \alpha) f(x|\alpha) + z(x, \alpha) \frac{\partial}{\partial \alpha} f(x|\alpha) = z(x, \alpha) f(x|\alpha)$$

$$= \int z'(x, \alpha) f(x|\alpha) dn + \int z''(x, \alpha) f(x|\alpha) dn$$

$$= \mathbb{E}(z'(x, \alpha)) + \mathbb{E}(z''(x, \alpha))$$

$$= \mathbb{E}(z'(x, \alpha)) + \text{Var}(z(x, \alpha)) \text{ as } \mathbb{E}(z(x, \alpha))$$

$$\text{Hence, } \text{Var}(z(x, \alpha)) = - \mathbb{E}(z'(x, \alpha))$$

\* Proof sketch : Assuming  $f(n|\theta)$  is the PDF

To see why  $\hat{\theta}_{MLE}$  is consistent, we note that  $\hat{\theta}$  is the value at  $\theta$  that maximizes :

$$\frac{1}{n} \ell(\theta) = \frac{1}{n} \log \left( \prod_{i=1}^n f(n|\theta) \right) = \frac{1}{n} \sum_{i=1}^n \log (f(n|\theta))$$

Suppose  $x_1, \dots, x_n \stackrel{iid}{\sim} f(n|\theta_0)$ . By LLN

$$\frac{1}{n} \sum_{i=1}^n \log (f(n|\theta)) \xrightarrow{\text{P}} \mathbb{E} (\log (f(n|\theta)))$$

$$\mathbb{E} (\log (f(n|\theta))) - \mathbb{E} (\log (f(n|\theta_0))) = \mathbb{E} \left( \log \left( \frac{f(n|\theta)}{f(n|\theta_0)} \right) \right)$$

By Jensen's Inequality ( $u \mapsto \log$  is concave):

$$\begin{aligned} \mathbb{E} \left( \log \left( \frac{f(n|\theta)}{f(n|\theta_0)} \right) \right) &\leq \log \left( \mathbb{E} \left( \frac{f(n|\theta)}{f(n|\theta_0)} \right) \right) \\ &= \log \left( \int \frac{f(n|\theta)}{f(n|\theta_0)} f(n|\theta_0) d\theta \right) \\ &= \log \left( \int f(n|\theta) d\theta \right) \\ &= \log (1) \\ &= 0 \end{aligned}$$

Hence,  $\ell$  is maximized at  $\theta_0$

By the definition of MLE :  $\ell'(\hat{\theta}) = 0$

Consistency at  $\hat{\theta}$  ensures that  $\hat{\theta}$  is close to the true value at  $\theta$ , i.e.  $\hat{\theta}_0$  with high probability.

$$\ell'(\theta_0) + (\hat{\theta} - \theta_0) \ell''(\theta_0) = 0 \quad \text{by Taylor's Expansion}$$

$$\hat{\theta} - \theta_0 = -\frac{\ell'(\theta_0)}{\ell''(\theta_0)}$$

$$\sqrt{n}(\hat{\theta} - \theta) = -\frac{\sqrt{n}\ell'(\theta_0)}{\sqrt{n}\ell''(\theta_0)}$$

$$\begin{aligned} \frac{1}{n} \ell''(\theta) &= \frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log(f(x_i|\theta)) \\ &= \frac{1}{n} \sum_{i=1}^n z^*(n, \theta) \end{aligned}$$

$$\xrightarrow{\text{IP}} E(z^*(X, \theta)) = -I(\theta)$$

$$\begin{aligned} \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial}{\partial \theta} \log(f(x_i|\theta)) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n z(n, \theta) \\ &\xrightarrow{} N(0, I(\theta)) \end{aligned}$$

By Slutsky's Lemma:  $\sqrt{n}(\hat{\theta} - \theta) \sim \frac{1}{I(\theta)} N(0, I(\theta))$

$$\sqrt{n}(\hat{\theta} - \theta) \sim N\left(0, \frac{1}{I(\theta)}\right)$$

\* Remark: 1)  $\hat{\theta}$  is asymptotically unbiased. More precisely the bias at  $\hat{\theta}$  is less than the order  $\frac{1}{\sqrt{n}}$ . Otherwise  $\sqrt{n}(\hat{\theta} - \theta_0)$  does not converge to a zero mean distribution.

2) The variance of  $\hat{\alpha}$  is approximately  $\frac{1}{n I(\alpha)}$ .  
 In particular, the standard error is of order  $\frac{1}{\sqrt{n}}$  and the variance is the main contributor to the  $MSE = \text{Var} + \text{Bias}^2$  of  $\hat{\alpha}$ .

3) If the true parameter is  $\alpha_0$ , then the sampling distribution of  $\hat{\alpha}$  is approximately  $N(\alpha_0, \frac{1}{n I(\alpha)})$

### \* Example : Sampling Distribution of Gamma

$$x_1, \dots, x_n \stackrel{iid}{\sim} \text{Gamma}(\alpha, 1)$$

$$\begin{aligned}\text{log}(f(n|\alpha)) &= \text{log}\left(\frac{1}{\Gamma(\alpha)} \cdot n^{\alpha-1} e^{-n}\right) \\ &= (\alpha - 1) \text{log } n - n - \text{log}(\Gamma(\alpha))\end{aligned}$$

$$z(n, \alpha) = \text{log}(n) - \frac{\Gamma'(\alpha)}{\Gamma(\alpha)}$$

$$z'(n, \alpha) = -\Psi'(\alpha), \quad \Psi(\alpha) = \frac{\Gamma'(\alpha)}{\Gamma(\alpha)}$$

$$I(\alpha) = -\mathbb{E}(z'(X, \alpha)) = \Psi'(\alpha)$$

$$\sqrt{n}(\hat{\alpha}_{MLE} - \alpha) \sim N\left(0, \frac{1}{\Psi'(\alpha)}\right)$$

$$\hat{\alpha}_{MLE} \sim N\left(\alpha, \frac{1}{n \Psi'(\alpha)}\right)$$

### \* Fisher Information for more than one parameter :

$$\text{Let } x_1, x_2, \dots, x_n \stackrel{iid}{\sim} f(y|\theta) \quad \{f(y|\theta) \mid \theta \in \mathbb{R}^k\} \text{ with } k$$

parameters and  $\hat{\alpha}_n$  is the MLE. We define F. I matrix  $I(\alpha) \in \mathbb{R}^{k \times k}$  is given by:

$$I_{(i,j)} = \text{Cov}\left(\frac{\partial}{\partial \alpha_i} \log(f(x|\alpha)), \frac{\partial}{\partial \alpha_j} \log(f(x|\alpha))\right)$$

$$= -\mathbb{E}\left(\frac{\partial^2}{\partial \alpha_i \partial \alpha_j} \log(f(x|\alpha))\right)$$

$\sqrt{n}(\hat{\alpha}_n - \alpha) \sim MVN(0, I(\alpha)^{-1})$  where  $I(\alpha)^{-1}$  is the  $k \times k$  matrix inverse of  $I(\alpha)$ .

### \* Example: Sampling Distribution of Gamma

$x_1, \dots, x_n \stackrel{iid}{\sim} \text{Gamma}(\alpha, \beta)$

$$\begin{aligned} \log(f(u|\alpha, \beta)) &= \log\left(\frac{\beta^\alpha}{\Gamma(\alpha)} \cdot u^{\alpha-1} e^{-\beta u}\right) \\ &= \alpha \log(\beta) + (\alpha-1) \log u - \beta u - \log(\Gamma(\alpha)) \\ &= S \end{aligned}$$

$$\frac{\partial^2}{\partial \alpha^2} S = -\psi'(\alpha)$$

$$\frac{\partial^2}{\partial \beta^2} S = \frac{\partial}{\partial \beta} \left( \frac{\alpha}{\beta} - u \right) = -\frac{\alpha}{\beta^2}$$

$$\frac{\partial^2 S}{\partial \alpha \partial \beta} = \frac{\partial}{\partial \alpha} \left( \frac{\alpha}{\beta} - u \right) = \frac{1}{\beta}$$

$$\text{Hence, } I(\alpha, \beta) = \begin{pmatrix} \psi'(\alpha) & -\frac{1}{\beta} \\ -\frac{1}{\beta} & \frac{1}{\beta^2} \end{pmatrix}$$

$$I'(\alpha, \beta) = \frac{\beta^2}{\alpha \Psi'(\alpha) - 1} \begin{pmatrix} \alpha/\beta^2 & \gamma_\beta \\ \gamma_\beta & \Psi'(\alpha) \end{pmatrix}$$

$$(\hat{\alpha}, \hat{\beta}) \sim \text{BVN}\left((\alpha, \beta), \frac{1}{n} I'(\alpha, \beta)\right)$$

$$\hat{\alpha} \sim N\left(\alpha, \frac{\beta^2}{\alpha \Psi'(\alpha) - 1} \cdot \frac{\alpha}{n \beta^2}\right)$$

$$\hat{\beta} \sim N\left(\beta, \frac{\alpha}{n(\alpha \Psi'(\alpha) - 1)}\right)$$

Hence, larger model has higher variability.

**Implication:** A complex model may better capture the true distribution of the data but they will be increasing the difficulty in estimating the parameters than those in a simple model.

**Q<sub>1</sub>:** Why is Fisher information  $I(\theta)$  called information?

**Q<sub>2</sub>:** Why should we use MLE for estimating  $\theta$ ?

Let  $X_1, \dots, X_n \stackrel{iid}{\sim} f(n|\theta_0)$ ,  $\ell(\theta) = \sum_{i=1}^n \log(f(x_i|\theta))$  and  $I(\theta_0) = -\mathbb{E}\left(\frac{\partial^2}{\partial \theta^2} \log(f(x_i|\theta))|_{\theta=\theta_0}\right) = -\frac{1}{n} \mathbb{E}(\ell''(\theta_0))$ . Fisher information matrix measures the expected curvature of the log likelihood function  $\ell(\theta)$  around the true parameter  $\theta = \theta_0$ .

**A<sub>1</sub>:** F. I quantifies the amount of information that each observation  $x_i$  contains about the unknown parameter.

## \* Theorem : Cramer Rao Lower Bound (CRLB)

Consider a parametric model  $\{f(n|\theta) \mid \theta \in \mathcal{S}\}$  satisfying the regularity conditions where  $\theta \in \mathbb{R}$ . Let  $T$  be an unbiased estimator at  $\theta$  based on the data  $x_1, \dots, x_n \stackrel{iid}{\sim} f(n|\theta)$ . Then,  $\text{Var}(T) \geq \frac{1}{n I(\theta)}$

A<sub>2</sub>: The CRLB ensures that no unbiased estimator can achieve asymptotically lower variance than the MLE. In general, MOM estimator have higher MSE than MLE for large  $n$ .

\* Efficiency : Let  $T_1$  and  $T_2$  be two unbiased estimators at  $\theta$ .

$$\text{Relative efficiency} = \frac{\text{Var}(T_1)}{\text{Var}(T_2)}.$$

An unbiased estimator is efficient if its variance is  $= \frac{1}{n I(\theta)}$ . We say MLE is asymptotically efficient.

\* MLE under model misspecification :

If the (parametric) model is incorrect, does MLE estimate  $\hat{\theta}$  still remain meaningful?

We have been measuring the error of an estimator  $\hat{\theta}$  by bias, variance and MSE.

\* Def : KL Divergence

If  $x_1, \dots, x_n \stackrel{iid}{\sim} g$  that is not in the model (parametric family  $\{f(n|\theta) \mid \theta \in \mathcal{S}\}$ ) then there is no true parameter  $\theta$  associated to  $g$ .

Now, we introduce a measure of "distance" between two PDFs.

For two PDFs  $f$  and  $g$ , the KL divergence from  $f$  to  $g$  is given by:

$$D_{KL}(g||f) = \int g(n) \log\left(\frac{g(n)}{f(n)}\right) dn.$$

Equivalently, if  $X \sim g$ , then  $D_{KL}(g||f) = E\left(\log\left(\frac{g(n)}{f(n)}\right)\right)$

1) If  $f = g$  then  $D_{KL}(g||f) = 0$

$$\text{2) } E\left(\log\left(\frac{g(x)}{f(x)}\right)\right) = E\left(-\log\left(\frac{f(x)}{g(x)}\right)\right)$$

$$\geq -\log\left(E\left(\frac{f(x)}{g(x)}\right)\right)$$

$$\geq -\log\left(\int \frac{f(n)}{g(n)} \cdot g(n) dn\right)$$

$$\geq -\log(1)$$

$$\geq 0$$

3) If  $f$  and  $g$  are constants then  $D_{KL}(g||f) = 0$

\* Example :

Suppose  $f = N(\mu_0, \sigma^2)$  and  $g = N(\mu_1, \sigma^2)$ . Calculating  $D_{KL}(g||f)$  if  $X \sim g$  and  $D_{KL}(f||g)$ .

$$D_{KL}(g||f) = E\left(\log\left(\frac{g(x)}{f(x)}\right)\right)$$

$$\log\left(\frac{g(n)}{f(n)}\right) = \log\left(\frac{e^{-\frac{(n-\mu_1)^2}{2\sigma^2}}}{e^{-\frac{(n-\mu_0)^2}{2\sigma^2}}}\right) = \frac{-(n-\mu_1)^2}{2\sigma^2} + \frac{(n-\mu_0)^2}{2\sigma^2}$$

$$= \frac{-n^2 + 2n\mu_1 - \mu_1^2 + n^2 - 2n\mu_0 + \mu_0^2}{2\sigma^2}$$

$$= \frac{(\mu_0^2 - \mu_1^2) + 2n(\mu_1 - \mu_0)}{2\sigma^2}$$

$$\begin{aligned} \mathbb{E} \left( \log \left( \frac{f(x)}{g(x)} \right) \right) &= \mathbb{E} \left( \frac{(\mu_0^2 - \mu_1^2) + 2n(\mu_1 - \mu_0)}{2\sigma^2} \right) \\ &= \frac{(\mu_0^2 - \mu_1^2)}{2\sigma^2} + \frac{(\mu_1 - \mu_0)}{\sigma^2} \cdot \mu_1 \\ &= \frac{\mu_0^2 - \mu_1^2 + 2\mu_1^2 - 2\mu_0\mu_1}{2\sigma^2} \\ &= \frac{\mu_0^2 - 2\mu_0\mu_1 + \mu_1^2}{2\sigma^2} \\ &= \frac{(\mu_0 - \mu_1)^2}{2\sigma^2} \end{aligned}$$

$$D_{KL}(f || g) = \mathbb{E} \left( \log \left( \frac{f(x)}{g(x)} \right) \right)$$

$$\begin{aligned} \log \left( \frac{f(n)}{g(n)} \right) &= \log \left( \frac{e^{-\frac{(n-\mu_0)^2}{2\sigma^2}}}{e^{-\frac{(n-\mu_1)^2}{2\sigma^2}}} \right) = \frac{-(n-\mu_0)^2}{2\sigma^2} + \frac{(n-\mu_1)^2}{2\sigma^2} \\ &= \frac{-n^2 + 2n\mu_0 - \mu_0^2 + n^2 - 2n\mu_1 + \mu_1^2}{2\sigma^2} \\ &= \frac{(\mu_1^2 - \mu_0^2) + 2n(\mu_0 - \mu_1)}{2\sigma^2} \end{aligned}$$

$$\begin{aligned} \mathbb{E} \left( \log \left( \frac{f(x)}{g(x)} \right) \right) &= \mathbb{E} \left( \frac{(\mu_1^2 - \mu_0^2) + 2n(\mu_0 - \mu_1)}{2\sigma^2} \right) \\ &= \frac{(\mu_1^2 - \mu_0^2)}{2\sigma^2} + \frac{(\mu_0 - \mu_1)}{\sigma^2} \cdot \mu_0 \\ &= \frac{\mu_1^2 - \mu_0^2 + 2\mu_0^2 - 2\mu_0\mu_1}{2\sigma^2} \\ &= \frac{\mu_0^2 - 2\mu_0\mu_1 + \mu_1^2}{2\sigma^2} \\ &= \frac{(\mu_0 - \mu_1)^2}{2\sigma^2} \end{aligned}$$

- \* Let  $x_1, x_2, \dots, x_n \sim g$ , we suppose that  $D_{KL}(g||f)$  has a unique minimum  $\theta = \theta^*$ . Then, under suitable regularity conditions, the MLE  $\hat{\theta}$  converges to  $\theta^*$  in probability as  $n \rightarrow +\infty$ . Then, we call  $f(n|\theta^*)$  as the KL projection of  $g$  onto the parametric model  $\{f(n|\theta) \mid \theta \in \mathbb{R}\}$ . In other words, the MLE is estimating the distribution in our model that is closest w.r.t KL divergence to  $g$ .

### \* Sandwich estimator of variance :

Let  $x_1, \dots, x_n \stackrel{iid}{\sim} g$ , how close is the MLE  $\hat{\theta}$  to this KL-projection.

$\theta = l'(\hat{\theta})$ , we write Taylor expansion around  $\hat{\theta} = \theta^*$

$$\theta \approx l'(\theta^*) + (\hat{\theta} - \theta^*) l''(\theta^*)$$

$$\Rightarrow (\hat{\theta} - \theta^*) = - \frac{l'(\theta^*)}{l''(\theta^*)}$$

$$\Rightarrow \sqrt{n}(\hat{\theta} - \theta^*) = \frac{-\sqrt{n}l'(\theta^*)}{\sqrt{n}l''(\theta^*)} \xrightarrow{\text{normalizing } z(x, \theta)} N(0, \text{Var}(z(x, \theta^*)))$$

$$= \frac{\sqrt{n}l'(\theta^*)}{-\mathbb{E}(z'(x, \theta^*))} \xrightarrow{} N\left(0, \frac{\text{Var}(z(x, \theta^*))}{\mathbb{E}(z'(x, \theta^*))^2}\right)$$

### \* Example :

$x_1, \dots, x_n \stackrel{iid}{\sim} \text{Exp}(\lambda)$

$$\theta = \frac{\partial}{\partial \lambda} \left( \log(\lambda^n e^{-\lambda \sum_{i=1}^n x_i}) \right) = \frac{\partial}{\partial \lambda} \left( n \log \lambda - \lambda \sum_{i=1}^n x_i \right) = \frac{n}{\lambda} - \sum_{i=1}^n x_i$$

$$\hat{\lambda}_{MLE} = \frac{1}{\bar{x}}$$

$$z(x, \lambda) = \frac{\partial}{\partial \lambda} (\log(f(n|\lambda))) = \frac{\partial}{\partial \lambda} (\log \lambda - \lambda n) = \frac{1}{\lambda} - n$$

$$z'(x, \lambda) = \frac{\partial}{\partial \lambda} z(x, \lambda) = \frac{-1}{\lambda^2}$$

$$\bar{z}(n, \lambda) = \frac{1}{\lambda} - \frac{1}{n} \sum_{i=1}^n u_i = \frac{1}{\lambda} - \bar{u}$$

$$\mathbb{E}(z(x, \alpha^*)) = \frac{1}{n} \sum_{i=1}^n \left( \frac{-1}{\lambda^2} \right) = \frac{-1}{\lambda^2} = -\bar{x}^2$$

$$\begin{aligned} \text{Var}(z(x, \alpha^*)) &= \frac{1}{n-1} \sum_{i=1}^n (z(x_i, \alpha^*) - \bar{z}(x_i, \alpha^*))^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n \left( \frac{1}{\lambda} - u_i - \frac{1}{\lambda} + \bar{u} \right)^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n (u_i - \bar{u})^2 \\ &= S_x^2 \end{aligned}$$

$$\text{Sandwich estimate of Variance : } \frac{\text{Var}(z(x, \alpha^*))}{\mathbb{E}(z(x, \alpha^*))^2} = \frac{S_x^2}{\bar{x}^4}$$

### \* **Plugin Estimator :**

So far, we know how to estimate  $\theta$ . Then, we learnt how to compute MLE. Then, we learnt the sampling distribution and variance for MLE. No unbiased estimator achieves smaller than that of the MLE for large  $n$  (CRLB).

Now, we are interested in estimating  $g(\theta)$  by  $g(\hat{\theta})$ , where  $\hat{\theta}$  (say MLE) is an estimate of  $\theta$ . This is called plugin estimator.

### \* **Example : Pareto Distribution**

$$f(u | \lambda, n) = \lambda^{-n} u^{-\lambda} n^{-\lambda-1} \Gamma(n/\lambda)$$



$(n_i, m_i) = \alpha \cdot n^{\alpha} \cdot m$ , where  $\alpha > 1$

Used as an income distribution and  $m$  represents the minimum possible income.

Barabasi Network Talk

One parameter power distribution :  $f(n) = \alpha \cdot n^{-\alpha-1} \quad (n \geq 1)$

$$\begin{aligned}
 E(X) &= \int_{\mathbb{R}} n f(n) \alpha \, dn \\
 &= \int_1^{+\infty} \alpha \cdot n^{-\alpha} \, dn \\
 &= \alpha \lim_{n \rightarrow +\infty} \left[ \frac{n^{-\alpha+1}}{-\alpha+1} \right]_1^n \\
 &= \alpha \lim_{n \rightarrow +\infty} \left( \frac{n^{1-\alpha}}{1-\alpha} + \frac{1}{\alpha-1} \right) \quad 1-\alpha < 0 \text{ as } \alpha > 1 \\
 &= \frac{\alpha}{\alpha-1}
 \end{aligned}$$

$$l(\alpha) = n \log(\alpha) - (\alpha+1) \sum_{i=1}^n \log(n_i)$$

$$\hat{\alpha} = \frac{n}{\sum_{i=1}^n \log(n_i)}$$

$$\hat{\alpha}_{MLE} = \frac{n}{\sum_{i=1}^n \log(n_i)}$$

Hence, the average income is estimated by :  $\frac{\hat{\alpha}_{MLE}}{\hat{\alpha}_{MLE}-1}$

We we don't have enough data but we know the distribution of the data so we can still estimate  $g(\alpha)$ .

\* Examples : Odds (Odds Ratio)

You play a game with a friend, where you flip a biased coin. If

If the coin lands heads, then you give your friend 1 AED. If the coin lands tails, your friend gives you  $n$  AED. What is the value of  $n$  that makes this a fair coin?

Fair game  $\Rightarrow$  earning = 0 AED

$$E(\text{winning}) = -P + (1-P)n = 0 \Rightarrow n = \frac{P}{1-P} \text{ odds of getting heads over tails}$$

$$X \sim \text{Ber}(P)$$

$$\ell(P) = \sum_{i=1}^n n_i \log(P) + \sum_{i=1}^n (1-n_i) \log(1-P)$$

$$\ell'(P) = \frac{\sum n_i}{P} - \frac{n - \sum n_i}{1-P} = 0$$

$$(1-P) \sum_{i=1}^n n_i - np + p \sum_{i=1}^n n_i = 0$$

$$\sum_{i=1}^n n_i - np = 0$$

$$\hat{P}_{MLE} = \bar{x}$$

Hence, the estimate  $\frac{\bar{x}}{1-\bar{x}}$

### \* Delta Method:

$$g(\hat{\theta}) = g(\theta_0) + (\hat{\theta} - \theta_0) g'(\theta_0) \quad \text{Taylor expansion around } \theta_0$$

$$g(\hat{\theta}) - g(\theta_0) = (\hat{\theta} - \theta_0) g'(\theta_0)$$

$$\Rightarrow \sqrt{n} (g(\hat{\theta}) - g(\theta_0)) = \sqrt{n} (\hat{\theta} - \theta_0) g'(\theta_0) \sim N\left(0, \frac{g'(\theta_0)^2}{I(\theta_0)}\right)$$

\* Using the delta method, we can quantify the uncertainty about the plugin estimator  $g(\hat{\theta})$  for large  $n$ .

\* Restrictions :  
 i)  $g : \mathbb{R} \rightarrow \mathbb{R}$  real valued  
 ii)  $g'(0_0) \neq 0$

$$\sqrt{n} (g(\hat{\theta}) - g(\theta_0)) \xrightarrow{\text{in distribution}} N(0, \frac{g'(\theta_0)^2}{I(\theta_0)}) \text{ as } n \rightarrow +\infty$$

\* Example : Log-odds

Let  $x_1, x_2, \dots, x_n \stackrel{\text{iid}}{\sim} \text{Ber}(p)$  and recall the estimate of log-odds  $\log\left(\frac{p}{1-p}\right)$  given by  $\log\left(\frac{\bar{x}}{1-\bar{x}}\right)$ .

By CLT :  $\sqrt{n}(\bar{x} - p) \sim N(0, p(1-p))$

$$(\log(p) - \log(1-p))' = \frac{1}{p} + \frac{1}{1-p} = \frac{1}{p(1-p)}$$

By the delta method :  $\sqrt{n}(\log\left(\frac{\bar{x}}{1-\bar{x}}\right) - \log\left(\frac{p}{1-p}\right)) \sim N(0, \frac{1}{p(1-p)})$

Interpretation : Our estimate log-odds of heads to tails is approximately normally distributed around the true odds  $\log\left(\frac{p}{1-p}\right)$  with variance  $\frac{1}{np(1-p)}$ .

Suppose we toss a coin 1000 times. We observe 600 heads.

$$\bar{x} = \hat{p} = \frac{600}{1000} = 0.6, \text{ we estimate log-odds by } \log\left(\frac{0.6}{0.4}\right) = 0.41$$

and standard error  $\sqrt{\frac{1}{n\bar{x}(1-\bar{x})}} \approx 0.2$ .

\* Example : Pareto Distribution

We showed that the average income is estimated by  $g(\theta) = \frac{\theta}{\theta-1}$   
 where  $\hat{\theta}_{MLE} = \frac{n}{m}$ . Also,  $g'(\theta) = -\frac{1}{(\theta-1)^2}$ .

$$\sum_{i=1}^n \text{log}(m_i)$$

$$(a-1)^2$$

Hence, by the delta method :  $\sqrt{n} \left( \frac{\hat{\alpha}}{\alpha-1} - \frac{\alpha}{\alpha-1} \right) \sim N(0, \frac{\alpha^2}{(\alpha-1)^4})$

$$\begin{aligned} I(\alpha) &= -E(\tilde{z}'(x, \alpha)) \\ &= -E\left(\frac{\partial^2}{\partial \alpha^2} (\text{log}(m) - (\alpha+1) \text{log}(n))\right) \\ &= -E\left(\frac{\partial}{\partial \alpha} \left(\frac{1}{\alpha} - \text{log}(n)\right)\right) \\ &= -E\left(-\frac{1}{\alpha^2}\right) \\ &= \frac{1}{\alpha^2} \end{aligned}$$

Hence,  $\sqrt{n}(\hat{\alpha} - \alpha) \sim N(0, \alpha^2)$ .

Finally,  $\frac{\hat{\alpha}}{\hat{\alpha}-1} \sim N\left(\frac{\alpha}{\alpha-1}, \frac{\alpha^2}{n(\alpha-1)^4}\right)$ .

### \* Theorem : General Cramér Rao Lower Bound

For a parametric model  $\{f(x|\alpha) \mid \alpha \in \mathcal{A}\}$  satisfying certain mild regularity conditions and  $g$  be any function differentiable on all of  $\mathcal{A}$  and  $T$  be any unbiased estimator of  $g(\alpha)$  based on data  $x_1, x_2, \dots, x_n \stackrel{iid}{\sim} f(x|\alpha)$ . Then,  $V(T) \geq \frac{g'(\alpha)^2}{n I(\alpha)}$ .

When an estimator  $T$  achieves this variance, it is called efficient.

### \* Properties of Estimators :

From what we learnt consistency and asymptotic efficiency are large sample property whereas unbiasedness and minimum variance are practical properties. Are there any more practical properties?

## \* Def : Minimum Variance Unbiased Estimator (MVUE)

Unbiased estimators can average fit the target for all the parameter values. This seems to be a reasonable constraint to impose on an estimator and indeed produces meaningful estimates in a wide variety of situations. In a large class of problems, it is possible to find an VE of  $g(\theta)$  that has the smallest possible variance among all these unbiased estimators. Such an estimate is called MVUE.

$$1) \mathbb{E}_\theta(S_n) = g(\theta)$$

$$2) \text{Var}_\theta(S_n) \leq \text{Var}_\theta(T_n) \text{ where } T_n \text{ is an unbiased estimator of } g(\theta)$$

## \* Example : Normal Distribution

$$x_1, \dots, x_n \sim N(\mu, \sigma^2)$$

$\bar{x}$ : sample mean

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\mathbb{E}(\bar{x}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(x_i) = \mu \quad \text{unbiased}$$

$$\text{Var}(\bar{x}) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(x_i) = \frac{\sigma^2}{n} \longrightarrow 0$$

Hence,  $\bar{x}$  is consistent & MVUE as it achieves the lowest variance.

$$\mathbb{E}(\hat{\sigma}^2) = \mathbb{E}\left(\frac{n-1}{n} s^2\right) = \frac{(n-1)\sigma^2}{n}$$

$$\mathbb{E}(s^2) = \frac{1}{n-1} \sum_{i=1}^n \mathbb{E}((x_i - \bar{x})^2) = \frac{1}{n-1} \sum_{i=1}^n \mathbb{E}((x_i - \mu + \mu - \bar{x})^2)$$

$$\begin{aligned}
&= \frac{1}{n-1} \left( \sum_{i=1}^n \mathbb{E}((x_i - \mu)^2) + \sum_{i=1}^n \mathbb{E}((\mu - \bar{x})^2) - 2 \sum_{i=1}^n \mathbb{E}((x_i - \mu)(\mu - \bar{x})) \right) \\
&= \frac{1}{n-1} \left( n \text{Var}(x_i) + n \text{Var}(\bar{x}) - 2n \mathbb{E}((\mu - \bar{x}^2)) \right) \\
&= \frac{1}{n-1} \left( \sigma^2(n-1) \right) \\
&= \sigma^2 \quad \text{MVUR}
\end{aligned}$$

Unbiased estimates may not always be better than biased estimators.

In ML, performance vs model complexity trade-off

\* Example :  $x_1, x_2, \dots, x_n \stackrel{iid}{\sim} \mathcal{U}(0, \theta)$ ,  $\theta > 0$   $\Omega = [0, +\infty]$ . A natural estimate of  $\theta$  is  $\hat{\theta}_{(n)} = \max\{x_i\}$

$\mathbb{E}(\bar{x}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(x_i) = \frac{\theta}{2}$  so  $\bar{x}$  is an unbiased estimator of  $\theta$ .

Show that  $\hat{\theta}_{(n)}$  outperforms  $\bar{x}$  by an order of magnitude in the sense of MLE. Then, find out the MVUE of  $\theta$ .

$$\mathbb{E}(\hat{\theta}_{(n)}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(x_i) = 0$$

$$\text{Var}(\hat{\theta}_{(n)}) = \frac{4}{n^2} \sum_{i=1}^n \text{Var}(x_i)$$

$$\begin{aligned}
&= \frac{4}{n} \text{Var}(x_i) \quad \mathbb{E}(x_i^2) = \int_0^\theta \frac{n^2}{\theta} du = \frac{n^3}{3\theta} \Big|_0^\theta = \frac{\theta^2}{3} \\
&= \frac{4}{n} \left( \frac{\theta^2}{3} - \frac{\theta^2}{4} \right)
\end{aligned}$$

$$= \frac{\alpha^2}{3n}$$

- $MSE(\bar{x}, \alpha) = \text{Var}(\bar{x}) + \text{Bias}^2 = \frac{\alpha^2}{3n}$

- $F_n(n) = P(X_{(n)} \leq n)$   
 $= (F_x(n))^n$

$$f_n(n) = n (F_x(n))^{n-1} \cdot f_x(n)$$

$$= n \left( \frac{n}{\alpha} \right)^{n-1} \cdot \frac{1}{\alpha} \mathbb{1}_{[0, \alpha]}$$

$$f_n(n) = \frac{n n^{n-1}}{\alpha^n} \mathbb{1}_{[0, \alpha]}$$

- $E(X_{(n)}) = \frac{n}{\alpha^n} \int_0^\alpha u^n du \quad E(X_{(n)}^2) = \frac{n}{\alpha^n} \int_0^\alpha u^{n+1} du$

$$= \frac{n}{\alpha^n} \left( \frac{u^{n+1}}{n+1} \Big|_0^\alpha \right) \quad = \frac{n}{\alpha^n} \left( \frac{u^{n+2}}{n+2} \Big|_0^\alpha \right)$$

$$= \frac{n \alpha}{n+1} \quad = \frac{n \alpha^2}{n+2}$$

- $\text{Var}(X_{(n)}) = \frac{n \alpha^2}{n+2} - \frac{n^2 \alpha^2}{(n+1)^2} = \frac{n \alpha^2}{(n+2)(n+1)^2}$

- $MSE(X_{(n)}, \alpha) = \frac{n \alpha^2}{(n+2)(n+1)^2} + \left( \frac{n \alpha}{n+1} - \alpha \right)^2$

$$= \frac{n^2 \alpha^2}{n(n+2)(n+1)^2} + \frac{\alpha^2}{(n+1)^2}$$

- $MSE\left((1 + \frac{1}{n}) X_{(n)}, \alpha\right) = \frac{(n+1)^2}{n^2} \left( \frac{n^2 \alpha^2}{n(n+2)(n+1)^2} \right) = \frac{\alpha^2}{n(n+2)}$

Because  $E\left((1 + \frac{1}{n}) X_{(n)}\right) = \frac{n+1}{n} \left( \frac{n \alpha}{n+1} \right) = \alpha$  so it is unbiased.

## \* Sufficient Statistic :

A: Given  $x_1, \dots, x_n$



B: Given  $T = \phi(x_1, \dots, x_n)$



Generally, A will be able to find a better estimator than B. However, if B will be able to do just as good as A where  $T = \phi(x_1, \dots, x_n)$  will in some sense summarize all the information contained in the random sample about  $\theta$  then knowledge of the whole data  $x_1, \dots, x_n$  will be irrelevant in the search for a good estimator of  $\theta$ . A statistic  $T$  with such property is called a sufficient statistic.

## \* Def: Sufficient Statistic

Let  $x_1, \dots, x_n$  be a random sample from a distribution indexed by a parameter  $\theta \in \Omega$ . Let  $T$  be a statistic. Suppose that, for every  $\theta \in \Omega$  and every possible value of  $T$ , the conditional distribution of  $x_1, \dots, x_n$  given  $T = t$  (at  $\theta$ ) depends only on  $t$  but not  $\theta$ .

In conclusion,  $T$  is sufficient for obtaining as much info about  $\theta$  as one could get from  $(x_1, x_2, \dots, x_n)$ .

If  $T$  is complete sufficient statistic, then  $E(h(T)) = 0$

Likelihood Theorem : Achieves the minimum variance  $\Rightarrow$  UMVUE

## \* Example :

$$x_1, \dots, x_n \sim \mathcal{P}(\lambda), \lambda > 0 \quad T = \sum_{i=1}^n x_i$$

$$\begin{aligned}
 P(\underline{x} = \underline{n} | T(x) = t) &= \frac{P(x = \underline{n}, T(x) = t)}{P(T(x) = t)} \\
 &= \frac{P(x = \underline{n})}{P(T(x) = t)} \quad \text{if } T(\underline{n}) = t \\
 &= \frac{e^{-n\lambda} \lambda^{\underline{n}}}{\pi(n_1!) \dots \pi(n_t!)} \cdot \frac{t!}{e^{-nt} \cdot (n\lambda)^t} \\
 &= \frac{t!}{\prod_{i=1}^n n_i! n^t} \quad \text{independent of } \lambda
 \end{aligned}$$

Hence, it is a sufficient statistic

Trivial sufficient statistic :  $(x_1, x_2, \dots, x_n)$   
 $(x_{(1)}, x_{(2)}, \dots, x_{(n)})$

\* Finding a sufficient Statistic :

\* Theorem : Neyman Fisher Factorization Theorem (NFFT)

Let  $x_1, x_2, \dots, x_n \stackrel{iid}{\sim} f(n, \theta)$  where  $\theta$  is unknown and  $\theta \in \Omega$ . A statistic  $T(\underline{x})$  is a sufficient statistic for  $\theta$  iff the joint PDF / PMF  $f_n(\underline{x}, \theta)$  of  $(x_1, \dots, x_n)$  can be factorized as follows :

$$f_n(n, \theta) = u(\underline{n}) v(T(\underline{x}), \theta)$$

where  $u$  is independent of  $\theta$  and  $v$  depends on  $\theta$  and  $n$  through  $T(\underline{n})$ . Both  $u$  and  $v$  are non-negative.

\* Example :

$x_1, \dots, x_n \sim f(\alpha)$

$$\begin{aligned}
 f_n(n, \alpha) &= \prod_{i=1}^n \frac{e^{-\alpha} \cdot \alpha^{x_i}}{x_i!} \\
 &= \frac{e^{-n\alpha} \cdot \alpha^{\sum x_i}}{\left(\prod_{i=1}^n x_i!\right)} \\
 &= \frac{1}{\prod_{i=1}^n x_i!} (e^{-n\alpha} \cdot \alpha^{\sum x_i}) \\
 &= u(n) \cdot v(\alpha, T(n)) \text{ where } T(n) = \sum_{i=1}^n x_i
 \end{aligned}$$

\* Example :

$$f_n(n|\alpha) = \left\{ \Gamma(\alpha)^n \cdot \left( \prod_{i=1}^n x_i \right)^\alpha \right\} \left\{ \beta^n \cdot e^{-\beta \sum x_i} \right\}$$

$u(n) \qquad \qquad \qquad v(\beta, \alpha, T)$

where  $T(n) = \sum_{i=1}^n x_i$  where  $\alpha$  is known and  $\beta$  unknown or  
 $T(n) = \prod_{i=1}^n x_i$  where  $\alpha$  is unknown and  $\beta$  known

\* Theorem : Rao Blackwell Theorem

For every  $\delta \in \Omega$ ,  $MSE(\delta_0(T), g(\alpha)) \leq MSE(\delta(x), g(\alpha))$ .

In summary Rao Blackwell Theorem shows how to improve upon an estimator that is not a function of a sufficient statistic by using a new estimator that is function of a sufficient statistic.

If we know a function  $S_\alpha(T) = E_{\alpha}(\delta(\tilde{x})|T(\tilde{x}))$  where  $T(\tilde{x})$  is a sufficient statistic and  $\delta$  is an estimator of  $g(\alpha)$ . The function  $S$  is indeed an estimator of  $g(\alpha)$  because it only depends on the observations  $\tilde{x}$  but not on  $\alpha$ .

## \* Confidence Intervals :

In a parametric model, let  $g(\theta)$  be a quantity of interest. Informally, a confidence interval for  $g(\theta)$  is a random interval calculated from the data that contains this value  $g(\theta)$  with a specified probability. For example, 90% CI contains  $g(\theta)$  with prob 0.9. It means if we conduct 100 different 90% CIs for  $g(\theta)$  using 100 independent sets of data, then we expect about 90 of them to contain  $g(\theta)$ . "Idea of repeated sampling".

- \* Let  $x_1, x_2, \dots, x_n$  be a data sample. By a random interval whose lower and upper endpoints  $L(x_1, \dots, x_n)$  &  $U(x_1, \dots, x_n)$  are random functions of  $(x_1, \dots, x_n)$ . The interval  $[L(x_1, \dots, x_n), U(x_1, \dots, x_n)]$  is a  $(1-\alpha) 100\%$  confidence interval of  $g(\theta)$  where  $P_\theta$  denotes the probability under  $x_1, \dots, x_n \stackrel{iid}{\sim} f(n, \theta)$

$$P_\theta(L(x_1, \dots, x_n) \leq g(\theta) \leq U(x_1, \dots, x_n)) = 1-\alpha$$

## \* Example :

Let  $x_1, \dots, x_n \sim N(\mu, \sigma^2)$  where  $\mu$  and  $\sigma^2$  are unknown. Compute the CI of  $\mu$  based on  $\bar{x}$ .

$$\bar{x} \sim N(\mu, \frac{\sigma^2}{n}) \text{ where } \frac{(n-1)s^2}{\sigma^2} \sim \chi^2_{n-1} \quad \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

$$\frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \cdot \sqrt{\frac{\sigma^2}{s^2(n-1)}} \sim t_{n-1}$$

Let  $t_{\alpha/2,n-1}$  and  $-t_{\alpha/2,n-1}$  be the upper and lower  $\frac{\alpha}{2}$  point by symmetry  
this indicates

$$P_{\mu,\sigma^2} \left( -t_{\alpha/2,n-1} \leq \frac{\bar{X} - \mu}{S/\sqrt{n}} \leq t_{\alpha/2,n-1} \right) = 1 - \alpha$$

$$P_{\mu,\sigma^2} \left( \bar{X} - \frac{S}{\sqrt{n}} t_{\alpha/2,n-1} \leq \mu \leq \bar{X} + \frac{S}{\sqrt{n}} t_{\alpha/2,n-1} \right) = 1 - \alpha$$

### \* Definition : Confidence Set

A subset  $S$  is said to constitute a confidence set at confidence  $(1-\alpha)$  if  $P(S \in \Omega) \geq 1 - \alpha \quad \forall \Omega \in \mathcal{Q}$ .

### \* Methods of finding Confidence Interval :

Let  $\theta$  be a parameter and  $T$  be a statistic based on a random sample of size  $n$  for the population. Most often it is possible to find a function  $\Psi(T, \theta)$  whose distribution is independent of  $\theta$ . Then,  
 $P(\Psi_{1-\alpha/2} < \Psi(T, \theta) < \Psi_{\alpha/2}) = 1 - \alpha$  which can be reformulated as  $P(\theta_1(T) \leq \theta \leq \theta_2(T)) = 1 - \alpha$  and the observed interval  $[\theta_1(T), \theta_2(T)]$  is the CI with confidence level  $(1-\alpha)$ .

### \* Example :

- Find the CI of  $\sigma^2$ :

Let  $T(\sigma^2) = S^2$  and  $\frac{(n-1)S^2}{\sigma^2} \sim \chi^2_{n-1}$

$$P_{\frac{1-\alpha}{2},n-1} \left( \frac{(n-1)S^2}{\sigma^2} < \chi^2_{\alpha/2,n-1} \right) = 1 - \alpha$$

$$P\left(\frac{(n-1)s^2}{\chi^2_{\alpha/2, n-1}} < \sigma^2 < \frac{(n-1)s^2}{\chi^2_{1-\alpha/2, n-1}}\right) = 1 - \alpha$$

Hence, the CI of  $\sigma^2$  :  $\left[ \frac{(n-1)s^2}{\chi^2_{\alpha/2, n-1}}, \frac{(n-1)s^2}{\chi^2_{1-\alpha/2, n-1}} \right]$

- Find the CI of  $\mu, \sigma^2$ ;

Boolsen Inequality :  $P(A \cap B) \geq P(A) + P(B) - 1$

$$\text{let } \alpha = \alpha_1 + \alpha_2$$

$$\begin{aligned} & P\left(\bar{x} - \frac{s}{\sqrt{n}} t_{\alpha/2, n-1} < \mu < \bar{x} + \frac{s}{\sqrt{n}} t_{\alpha/2, n-1}\right) + P\left(\frac{(n-1)s^2}{\chi^2_{\alpha/2, n-1}} < \sigma^2 < \frac{(n-1)s^2}{\chi^2_{1-\alpha/2, n-1}}\right) - 1 \\ &= P(A) + P(B) - 1 \\ &= 1 - \alpha_1 + 1 - \alpha_2 - 1 \\ &= 1 - (\alpha_1 + \alpha_2) \\ &= 1 - \alpha \end{aligned}$$

By Boolsen inequality :  $P(A \cap B) \geq 1 - \alpha$

CI of  $\mu$  and  $\sigma^2$  :

$$P\left(\bar{x} - \frac{s}{\sqrt{n}} t_{\alpha/2, n-1} \leq \mu \leq \bar{x} + \frac{s}{\sqrt{n}} t_{\alpha/2, n-1}, \frac{(n-1)s^2}{\chi^2_{\alpha/2, n-1}} \leq \sigma^2 \leq \frac{(n-1)s^2}{\chi^2_{1-\alpha/2, n-1}}\right) = 1 - \alpha$$

\* Asymptotic Confidence Interval :

If  $n$  is large we can apply the CLT :  $\sqrt{n}(\bar{x} - \mu) \rightarrow N(0, \sigma^2)$  and  $s^2$  is a consistent estimator of  $\sigma^2$   $\frac{\sqrt{n}(\bar{x} - \mu)}{s} = \frac{\mu}{s} \cdot \frac{\sqrt{n}(\bar{x} - \mu)}{\sigma} \xrightarrow{n \rightarrow +\infty} N(0, 1)$

$$P\left(\bar{x} - \frac{s}{\sqrt{n}} z_{\alpha/2} < \mu < \bar{x} + \frac{s}{\sqrt{n}} z_{\alpha/2}\right) = 1 - \alpha. \text{ As } n \rightarrow +\infty t_{\alpha/2, n-1} \rightarrow z_{\alpha/2}.$$

### \* Example :

Let  $x_1, \dots, x_n$  find the asymptotic CI for  $\lambda$ .

$$T(x) = \bar{x} \Rightarrow E(\bar{x}) = \lambda \quad \text{Var}(\bar{x}) = \frac{\lambda}{n}$$

$$\text{By CLT: } \frac{\sqrt{n}(\bar{x} - \lambda)}{\sqrt{\lambda}} \sim N(0, 1) \text{ as } n \rightarrow +\infty$$

$$\text{Hence, } P(-z_{\alpha/2} \leq \frac{\sqrt{n}(\bar{x} - \lambda)}{\sqrt{\lambda}} \leq z_{\alpha/2}) = 1 - \alpha$$

$$\text{Therefore, the CI of } \lambda \text{ at level } (1 - \alpha) : \bar{x} \pm z_{\alpha/2} \sqrt{\frac{\lambda}{n}}$$

### \* General Asymptotic CI for MLE :

We know that  $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{} N(0, I(\theta)^{-1})$  as  $n \rightarrow +\infty$ . We estimate  $I(\theta)$  by  $I(\hat{\theta})$  by continuous mapping theorem  $I(\hat{\theta}) \xrightarrow{P} I(\theta)$ , this interval is called Wald interval.

$$\sqrt{n} I(\theta)(\hat{\theta} - \theta) = \boxed{\frac{\sqrt{I(\hat{\theta})}}{\sqrt{I(\theta)}}} \cdot \sqrt{n} I(\theta)(\hat{\theta} - \theta) \xrightarrow{n \rightarrow +\infty} N(0, 1)$$

↓  
converges to 1 in P

$$\Rightarrow P(-z_{\alpha/2} \leq \sqrt{n} I(\hat{\theta})(\theta - \hat{\theta}) \leq z_{\alpha/2}) = 1 - \alpha$$

$$\Rightarrow \text{CI} : \hat{\theta} \pm \frac{z_{\alpha/2}}{\sqrt{n} I(\hat{\theta})}$$

### \* Example :

$x_1, \dots, x_n \sim \text{Ber}(p)$  CL of  $g(p) = \log\left(\frac{p}{1-p}\right)$ .

- $\hat{p} = \bar{x} \Rightarrow g(\hat{p}) = \log\left(\frac{\bar{x}}{1-\bar{x}}\right)$

$$\prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}$$

- $g'(p) = \frac{(1-p)}{(1-p)^2 \cdot p} = \frac{1}{p(1-p)}$

By the delta method:  $\sqrt{n}(g(\hat{p}) - g(p)) \rightarrow N(0, \frac{g'(p)^2}{I(p)})$

$$\sqrt{n}(g(\hat{p}) - g(p)) \rightarrow N(0, \frac{1}{p(1-p)})$$

Hence, the CI of  $g(\hat{p})$ :  $g(p) \pm z_{\alpha/2} \sqrt{\frac{1}{n\hat{p}(1-\hat{p})}}$  by the continuous mapping theorem and Slutsky's lemma.

$$\sqrt{n\hat{p}(1-\hat{p})}(g(\hat{p}) - g(p)) \rightarrow N(0, 1)$$

$$\sqrt{\frac{n\hat{p}(1-\hat{p})}{n\hat{p}(1-\hat{p})}} \sqrt{n\hat{p}(1-\hat{p})}(g(\hat{p}) - g(p)) \rightarrow N(0, 1)$$

$$\Rightarrow \sqrt{n\hat{p}(1-\hat{p})}(g(\hat{p}) - g(p)) \rightarrow N(0, 1)$$

$$\Rightarrow P(-z_{\alpha/2} \leq \sqrt{n\hat{p}(1-\hat{p})}(g(\hat{p}) - g(p)) \leq z_{\alpha/2}) = 1-\alpha$$

### \* Problem:

Let  $x_1, x_2, \dots, x_n \stackrel{iid}{\sim} \text{Ber}(p)$  be  $n$  tosses of a biased coin and let  $\hat{p} = \bar{x}$ . In this problem, we explore two different ways to construct a 95% CI for  $p$ .

By CLT:  $\sqrt{n}(\hat{p} - p) \sim N(0, p(1-p))$

(1)

a) Use the plugin estimate  $\hat{p}(1-\hat{p})$  to obtain a 95% CI for  $p$ .

$$\frac{\sqrt{n}(\hat{p}-p)}{\sqrt{\hat{p}(1-\hat{p})}} \rightarrow N(0,1)$$

$$P(-z_{0.025} \leq \frac{\sqrt{n}(\hat{p}-p)}{\sqrt{\hat{p}(1-\hat{p})}} \leq z_{0.025}) = 0.95$$

$$P\left(\hat{p} - z_{0.025} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + z_{0.025} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right) = 0.95$$

$$p \in \left(\hat{p} \pm z_{0.025} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right)$$

b) Instead of using the plugin estimate  $\hat{p}(1-\hat{p})$ , note that (1) implies that for large  $n$ ,

$$P(-\sqrt{\hat{p}(1-\hat{p})} z_{\alpha/2} \leq \sqrt{n}(\hat{p}-p) \leq \sqrt{\hat{p}(1-\hat{p})} z_{\alpha/2}) \approx 1-\alpha$$

Solve this equation:

$$\sqrt{n}(\hat{p}-p) = \sqrt{p(1-p)} z_{\alpha/2} \text{ for } p \text{ in terms of } \hat{p}$$

$$\sqrt{n}(\hat{p}-p) = -\sqrt{p(1-p)} z_{\alpha/2} \text{ for } p \text{ in terms of } \hat{p}.$$

to obtain a different 95% CI for  $p$ .

$$n(\hat{p}-p)^2 \leq p(1-p) z_{\alpha/2}^2$$

$$\Rightarrow (n + z_{\alpha/2}^2)p^2 - (dn\hat{p} + z_{\alpha/2}^2)p + n\hat{p}^2 \leq 0, \text{ put } \alpha=0.05$$

$$p \in \left( \frac{\hat{p} + \frac{z_{\alpha/2}^2}{2n} \pm z_{\alpha/2} \sqrt{\hat{p} \frac{(1-\hat{p})}{n} + \frac{z_{\alpha/2}^2}{4n^2}}}{1 + \frac{z_{\alpha/2}^2}{n}} \right)$$

c) Perform a simulation study to determine the true coverages of the confidence intervals in points (a) and (b) for the 9 combinations of sample size  $n = 10, 40, 100$  and  $p = 0.1, 0.3, 0.5$  (for each perform at least  $B = 100,000$  simulations). In each simulation, you may estimate  $\hat{p}$  directly from  $\frac{1}{n} \text{Bin}(n, p)$ . Report the simulated coverage levels in two tables and answer which CIs for  $p$  yield the true coverage close to 95% for small values of  $n$ ?

## \* Bayesian Analysis :

So far we assumed  $\theta$  is unknown but non-random quantity - it is actually some fixed value describing the true distribution of data and our goal was to estimate  $\theta$ , "Frequentist paradigm of S.I".

Now, in Bayesian Inference, we model  $\theta$  as a random variable. The Bayesian Paradigm naturally incorporates our prior belief about the unknown parameter and update this belief based on observed data.

## \* Prior and Posterior Distributions :

Let  $X, Y$  be two RVs having PDF / PMF  $f_{x,y}(u,y)$ .

$$\text{Marginal: } f_x(u) = \begin{cases} \int f_{x,y}(u,y) dy, & \text{PDF} \\ \sum_{y \in S_y} f_{x,y}(u,y), & \text{PMF} \end{cases}$$

$$\text{Conditional: } f_{y|x}(y|u) = \frac{f_{x,y}(u,y)}{f_x(u)}$$

$$f(u,y) = f(y|u) \cdot f(u) = f(u|y) f(y)$$

In Bayesian, before the data is observed, the unknown parameter is modeled as a random variable having a probability distribution,  $f(\theta)$  called prior distribution. This distribution represents our belief about the value of this parameter.

Instead of  $\gamma$ , we have the two RVs  $X$  and  $\Theta$ :

$$f_{x,\theta}(n, \theta) = f_{x|\theta}(n, \theta) f_\theta(\theta)$$

$$\underset{\Theta|x}{f(\theta|n)} = \frac{f_{x|\theta}(n, \theta)}{f_x(n)}$$

$$\underset{\text{updated knowledge}}{=} \frac{f_{x|\theta}(n, \theta) f_\theta(\theta)}{\int f_{x|\theta}(n, \theta') f_\theta(\theta') d\theta'}$$

$$\underset{\text{of } \theta}{=} \frac{f_{x|\theta}(n, \theta) f_\theta(\theta)}{\int f_{x|\theta}(n, \theta') f_\theta(\theta') d\theta'} \rightarrow \text{Bayes Theorem}$$

$$\propto \underset{x|\theta}{f(n, \theta)} \cdot \underset{\theta}{f_\theta(\theta)} \rightarrow \text{Prior Knowledge of } \theta$$

Posterior  $\propto$  Likelihood. Prior

- Uniform  $[0, 1]$  (Uniform Prior) is similar to using a frequentist method.

\* Example :

Let  $P$  be  $[0, 1]$  be the probability of heads for a biased coin and let  $x_1, \dots, x_n$  be the outcome of  $n$  tosses of this coin. If we don't have any prior information about  $p$ , we might choose for its prior distribution of  $\text{Unif}[0, 1]$  having PDF  $f_p(p) = 1_{[0,1]}$ .

$$f_{x,p}(n_1, \dots, n_n | p) = f_{x|p}(n_1, \dots, n_n | p) \cdot f_p(p)$$

$$= \prod_{i=1}^n p^{n_i} (1-p)^{n_i} \mathbb{1}_{[0,1]}$$

$$= p^{\sum n_i} (1-p)^{n - \sum n_i} \mathbb{1}_{[0,1]}$$

$$f_x(n_1, \dots, n_n) = \int_0^1 p^s (1-p)^{n-s} dp, \quad s = \sum_{i=1}^n n_i$$

$$\sim \text{Beta}(s+1, n-s+1)$$

$$f_{p|x}(p | n_1, \dots, n_n) = \frac{p^s (1-p)^{n-s} \mathbb{1}_{[0,1]}}{B(s+1, n-s+1)},$$

$$B(s+1, n-s+1) = \frac{\Gamma(s+1)\Gamma(n-s+1)}{\Gamma(n+2)}$$

$$p | n_1, \dots, n_n \sim \text{Beta}(s+1, n-s+1)$$

For logistic regression, Beta prior is used as it is between [0,1]

### \* Example :

Suppose now, we have a prior belief that  $p$  is close to  $\frac{1}{2}$ . For example  $B(\alpha, \alpha)$  has the mean  $\frac{\alpha}{\alpha+\alpha} = \frac{1}{2}$  and variance  $\frac{1}{8\alpha+4}$ . The constant  $\alpha$  may be chosen depending on how confident we are, a priori, that  $P$  is  $\frac{1}{2}$  ( $\alpha = 1 \Rightarrow$  unit [0,1] prior,  $\alpha > 1 \Rightarrow$  prior is more concentrated around  $\frac{1}{2}$ ). Can you calculate the posterior distribution.

$$f_{x|p}(n_1, \dots, n_n | p) = \prod_{i=1}^n p^{n_i} (1-p)^{n_i} = p^s (1-p)^{n-s}$$

$$f_p(p) = \frac{1}{B(\alpha, \alpha)} \cdot p^{\alpha-1} \cdot (1-p)^{\alpha-1}$$

$$f_{p|x}(n_1, \dots, n_n | p) \propto p^{s+\alpha-1} \cdot (1-p)^{n+\alpha-s-1} \sim \text{Beta}(s+\alpha, n+\alpha-s)$$

If  $\alpha = 1$ , we go back to the other case.

\* **Conjugate Prior** : When the posterior has the same distributions as the prior (Ex. 2 Beta) then this type of prior is called conjugate prior.

\* **Example :**

Let  $\lambda \in [0, +\infty]$  be the parameter of poisson model  $x_1, \dots, x_n \stackrel{iid}{\sim} \text{Pois}(\lambda)$ . Assume the prior for  $\lambda$  is Gamma( $\alpha, \beta$ ). Show that Gamma is a conjugate prior for  $\lambda$ .

$$f_{x|n}(\lambda) = \prod_{i=1}^n \frac{e^{-\lambda} \cdot \lambda^{x_i}}{x_i!} = \frac{e^{-n\lambda} \cdot \lambda^{\sum x_i}}{\left( \prod_{i=1}^n x_i! \right)} \propto e^{-n\lambda} \cdot \lambda^s$$

$$f_\lambda(\lambda) = \frac{\beta^\alpha}{\Gamma(\alpha)} \cdot \lambda^{\alpha-1} \cdot e^{-\beta\lambda} \propto \lambda^{\alpha-1} \cdot e^{-\beta\lambda}$$

$$\underset{x|n}{f(\lambda|x)} \propto \lambda^{s+\alpha-1} \cdot e^{-(\beta+n)\lambda} \sim \text{Gamma}(\alpha+s, \beta+n)$$

Hence,  $\lambda$  is a conjugate prior.

\* **Point and Interval Estimation :**

In many practical applications, we would like to have a single estimate  $\hat{\theta}$  as well as interval describing uncertainty about  $\theta$ .

Posterior Mean/Mode can be used as the estimate of  $\hat{\theta}$ . For the interval  $P(\theta \in I | X) = 1 - \alpha$  will give us Bayesian Credible interval.

$$\hat{P} = \frac{n+\alpha}{n+2\alpha} \text{ is the point estimate for Ex. 2}$$

We can use 0.05 and 0.95 quantiles of Beta( $\alpha$ ,  $n-s+\alpha$ ) as 90% Bayesian credible interval for  $P$ .

Frequentist	Bayesian
Data are a repeatable random sample	Parameters are unknown and described probabilistically
Parameters are fixed	Analysis is done conditionally on the observed data i.e. data is treated as fixed
Underlying parameters remain constant during this repeatable process	

### \* Example :

Sampling from an exponential distribution. Suppose that the distribution of the lifetime of fluorescent tubes of a certain type is the exponential distribution with parameter  $\theta$ . Suppose that  $x_1, \dots, x_n$  is a random sample of lamps of this type. Assume that  $\theta \sim \text{Gamma}(\alpha, \beta)$  for known  $\alpha, \beta$ .

$$f(x_1, \dots, x_n | \theta) = \prod_{i=1}^n \theta e^{-\theta x_i} = \theta^n e^{-\theta \sum x_i}, \quad s = \sum_{i=1}^n x_i$$

$$f_\theta(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \cdot \theta^{\alpha-1} \cdot e^{-\beta\theta} \propto \theta^{\alpha-1} e^{-\beta\theta}$$

$$\alpha+n-1 \quad -(\beta+s)\theta$$

$$f(\theta | x_1, \dots, x_n) \propto e^{-\theta} \theta^{n-1}$$

Posterior  $\Theta | x = n \sim \text{Gamma}(\alpha + n, \beta + s)$

\* **Def** : Loss Function

A loss function is a real valued function of two variables  $L(\theta, a)$  where  $\theta \in \mathbb{R}$  and  $a \in \mathbb{R}$ .

Squared loss :  $L(\theta, a) = (\theta - a)^2$

Absolute loss :  $L(\theta, a) = |\theta - a|$

\* **Def** : Bayes Estimator

$T(x)$  is an estimator of  $\theta$  and  $f(\cdot)$  is a prior PDF / PMF of  $\theta \in \mathbb{R}$ .

Consider the problem of estimating  $\theta$  without being able to observe the data. If the statistician chooses a particular estimate  $a$ , then the expected loss will be

$$\mathbb{E}(L(\theta, a)) = \int_{\mathbb{R}} L(\theta, a) \cdot f(\theta) d\theta$$

The statistician wishes to choose an estimate  $a$  for which the expected loss is minimum.

If the statistician observes the data  $n$  then  $f(\cdot | n)$  denotes the posterior PDF / PMF. The expected loss becomes :

$$\mathbb{E}(L(\theta, a | n)) = \int_{\mathbb{R}} L(\theta, a) \cdot f(\theta | n) d\theta$$

For each possible value  $n$  of  $X$ . Let  $\delta(n)$  denote a value of the estimate

a for which the expected loss is minimum. Then, the function  $\delta(\underline{x})$  is called the Bayes estimator of  $\theta$ .

\* **Remark :** The Bayes estimator is an estimator chosen to minimize the posterior mean of some measure of how far the estimator is from the parameter.

\* **Corollary :** Property of Posterior Mean (squared)

Let  $\theta \in \Omega \subset \mathbb{R}$ . Suppose the squared loss function is defined and we know the posterior mean i.e.  $E(\theta | \underline{x})$  is finite. Then, the Bayes estimator of  $\theta$  is :  $\delta(\underline{x}) = E(\theta | \underline{x})$

\* **Example :**

The estimate of  $\hat{p} = \frac{\alpha + s}{n + 2\alpha}$  for the second example by this corollary.

\* **Corollary :** Property of Posterior Median (Absolute)

Let  $\theta \in \Omega \subset \mathbb{R}$ . Suppose the absolute loss function is defined and we know the posterior median is finite. Then, the Bayes estimator of  $\theta$  is  $\delta^*(\underline{x})$ .

\* **Example :**

Suppose  $x_1, \dots, x_n$  be a random sample from  $N(\theta, \sigma^2)$  where  $\theta$  is unknown. Prior for  $\theta \sim N(\mu_0, v_0)$ . Find out the posterior distribution of  $\theta | x_1 = n_1, \dots, x_n = n_n$  ?

$$f_{x|\theta}(n| \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(n_i - \theta)^2}{2\sigma^2}\right)$$

$$\propto \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (n_i - \theta)^2\right)$$

$$\propto \exp\left(-\frac{1}{2\sigma^2} \left( \sum_{i=1}^n (n_i - \bar{n}) + n(\theta - \bar{n})^2 \right)\right)$$

$$\propto \exp\left(-\frac{n}{2\sigma^2} (\theta - \bar{n})^2\right) \quad \text{because it is independent of } \theta$$

$$f_\theta(\theta) = \frac{1}{\sqrt{2\pi v_0}} \cdot \exp\left(-\frac{1}{2v_0} (\theta - \mu_0)^2\right) \propto \exp\left(-\frac{1}{2v_0} (\theta - \mu_0)^2\right)$$

$$f_{\theta|x}(n|\theta) \propto \exp\left(-\frac{n}{2\sigma^2} (\theta - \bar{n})^2 - \frac{1}{2v_0} (\theta - \mu_0)^2\right)$$

$$\propto \exp\left(-\frac{1}{2v_1} (\theta - \mu_1)^2\right)$$

Hence,  $\theta | x_1 = n_1, \dots, x_n = n_n \sim N(\mu_1, v_1)$ ,  $\mu_1$  is the Bayes estimator

$$v_1^{-2} = \frac{\sigma^2 v_0^{-2}}{\sigma^2 + nv_0^{-2}}$$

$$\mu_1 = \frac{\sigma^2 \mu_0 + n \theta_0 v_0^{-2} \bar{n}}{\sigma^2 + nv_0^{-2}}$$

\* Normal Approximation for large  $n$  :

Consider Bayesian inference applied with the prior  $f_\theta(\theta)$  for a parametric model  $f_{x|\theta}(n|\theta)$ .

Log likelihood :  $\ell(\theta) = \sum_{i=1}^n \log(f_{x|\theta}(n_i|\theta))$

$$\text{Posterior} : f_{\theta|x}( \theta | n_1, \dots, n_m) = \frac{f_\theta(\theta)}{f_\theta(\hat{\theta})} f_{x|\theta}(n_1, \dots, n_m | \theta) \\ = f_\theta(\theta) \cdot \exp(l(\theta))$$

Applying Second order Taylor expansion of  $l(\theta)$  around the MLE of  $\theta = \hat{\theta}$ :

$$l(\theta) \approx l(\hat{\theta}) + (\theta - \hat{\theta}) \overbrace{l'(\hat{\theta})}^{\theta} + \frac{(\theta - \hat{\theta})^2}{2} l''(\hat{\theta}) \\ \approx l(\hat{\theta}) - \frac{n}{2} (\theta - \hat{\theta})^2 I(\hat{\theta}) \quad \text{for large } n$$

$$f_{\theta|x}( \theta | n) \propto \exp \left( l(\hat{\theta}) - \frac{n}{2} (\theta - \hat{\theta})^2 I(\hat{\theta}) \right) f_\theta(\theta) \\ \downarrow \text{function of data, independent of } \theta \\ \propto \exp \left( -\frac{n}{2} (\theta - \hat{\theta})^2 I(\hat{\theta}) \right) f_\theta(\theta) \\ \propto \exp \left( -\frac{1}{2} \frac{(\theta - \hat{\theta})^2}{I(\hat{\theta})} \right)$$

$$\text{Hence, } f_{\theta|x}( \theta | n) \sim N \left( \hat{\theta}, \frac{1}{n I(\hat{\theta})} \right).$$

\* **Remark K:** To summarize the posterior mean for  $\theta$  is, for large  $n$ , is approximately the MLE  $\hat{\theta}$ . Furthermore, a  $100(1-\alpha)$  Bayesian Credible interval is approximately given by  $\hat{\theta} \pm z_{\alpha/2} \frac{1}{\sqrt{n I(\hat{\theta})}}$  which is exactly the  $100(1-\alpha)\%$  Wald CI for  $\theta$ . In this sense, frequentist & Bayesian method yield similar inference for large  $n$ .